

2nd week of Jan: pieces of the project

I. The beginning: Monge's Problem

Take two sets of points with the same number of elements.

$$(\text{MP}) \min_{\sigma \in \Sigma_n} \sum_{i=1}^n c(x_i, y_{\sigma(i)}) \quad (\text{cost of transporting points of } \bar{X} \text{ to points of } \bar{Y})$$

Spoiler: We found an algorithm in $O(n^3)$ - Hungarian, Auction (BentseKas)

1D case:  $c(x, y) = h(x-y)$, h convex, e.g. $c(x, y) = |x-y|^p$, $p \geq 1$

Prop.: For h convex, any optimal σ is monotonic (respects the sorting)
 \Rightarrow algorithm in $O(n \log(n))$

Moreover, if h is strictly convex and points are disjoint, then $\exists!$ OT σ .

RK: Problem with $c(x, y) = |x-y|$  \leftarrow infinite solutions

Prop.: $X = Y = \mathbb{R}^d$ (ambient space of points) $c(x, y) = \|x-y\|^p$
 $p=1 \Rightarrow$ No crossing allowed

lets generalise:

Polish space

X a complete metric space, separable: $\alpha \in \mathcal{M}(X)$, $\beta \in \mathcal{M}(Y)$ Borel measures = Radon measures

Borel Polish

Notation for discrete measures: $\alpha = \sum_{i=1}^m a_i \delta_{x_i}$, $\beta = \sum_{j=1}^n b_j \delta_{y_j}$

Push-forwards: $T: X \rightarrow Y$ some map, $T_\# : \mathcal{M}(X) \rightarrow \mathcal{M}(Y)$ linear

Dof.: $T_\#(\delta_x) = \delta_{T(x)}$, $T_\# \left(\sum_i a_i \delta_{x_i} \right) = \sum_i a_i \delta_{T(x_i)}$

Dof. 1: $T_\# \alpha(B) := \alpha(T^{-1}(B))$, $\forall B \subset Y$ measurable

Dof. 2: $\forall g \in \mathcal{B}_c(Y)$, $\int g(y) d\beta(y) := \int g(T(x)) d\alpha(x)$

Prop: if $T: X \rightarrow Y$ is a diffeo. and α has density ρ_α w.r.t. λ then β has a density ρ_β w.r.t. λ and

$$\rho_\beta(x) = |\det(DT(x))| \cdot \rho_\alpha(T(x)).$$

Now:

$$(\text{MP}): \min_{\substack{\alpha \in \mathcal{P}(X) \\ \beta \in \mathcal{P}(Y)}} \inf_{T_\# \alpha = \beta} \int c(x, T(x)) d\alpha(x)$$

Rq.: $T_\#$ can only diminish the size of the support.

Prop. / RK: if $\begin{cases} \alpha = \frac{1}{n} \sum_i \delta_{x_i} \\ \beta = \frac{1}{m} \sum_j \delta_{y_j} \end{cases}$, then $T_\# \alpha = \beta \Leftrightarrow T$ induces a permutation ($T(x_i) = y_{\sigma(i)}$)

Theo: (Brenier)

$X = \mathbb{R}^d$, $C(x, y) = \|x - y\|^2$ (we can take manifolds and $d_{\text{eucl}}(x, y)^p$, $p > 1$), α has a density w.r.t. Leb.

i) \exists OT map T^* ($\alpha \rightarrow \beta$ c.c.)

ii) T^* is the only $T = \nabla \psi$, ψ convex, $(\nabla \psi)_\# \alpha = \beta$.

RK: ψ convex, $T = \nabla \psi \Leftrightarrow T \nearrow$

RK: if $\nabla \psi(x) = T_x$ is linear, then ψ is quadratic

RK: ψ convex $\Rightarrow \psi$ is a.e. differentiable $\Rightarrow \nabla \psi$ exists a.e.

Ex: $A = D$, $c = \|x - y\|^2$

$\alpha, \beta \in P(\mathbb{R})$, α has density $C_\alpha > 0$. C_α is the c.d.f., $Q_\alpha = C_\alpha^{-1}$ the quantile function.

Prop: $(C_\alpha)_\# \alpha = U_{[0,1]}$, $(Q_\alpha)_\# U_{[0,1]} = \alpha$.

Show it

Theo: By Brenier, $Q_\beta \circ C_\alpha = T$ is the unique OT from α to β .

For ex., take $\alpha = \mathcal{N}(m_\alpha, G_\alpha^2)$, $\beta = \mathcal{N}(m_\beta, G_\beta^2)$

Prop: $T(x) = \frac{G_\alpha}{G_\beta} (x - m_\alpha) + m_\beta$, $T_\# \alpha = \beta$

Show it

Cor: T is the unique OT. $W_2(\alpha, \beta) = \left((m_\alpha - m_\beta)^2 + (G_\alpha^2 - G_\beta^2)^2 \right)^{\frac{1}{2}}$

RK: $W_p(\alpha, \beta) = \inf_{T_\# \alpha = \beta} \int \|x - T(x)\|^p d\alpha(x)$ (The Wasserstein (?) distance).

Brenier

Note: We can transform this into a PDE problem, but this isn't the right formulation for the problem we want to study.

Ex: Gaussians

$\alpha = \mathcal{N}(m_\alpha, \Sigma_\alpha)$, $\beta = \mathcal{N}(m_\beta, \Sigma_\beta)$, $\Sigma = U \text{diag}(G^2) U^\top$

Ansatz: $T(x) = \overset{?}{A}(x - m_\alpha) + m_\beta$

Lemma: $T_\# \alpha = \beta \Leftrightarrow A \Sigma_\alpha A^\top = \Sigma_\beta$

Proof: $T_\# \alpha = \mathcal{N}(m, \Sigma)$, $m = \mathbb{E}[Tx]$, $X \sim \alpha$, $TX \sim T_\# \alpha$

$$m = m_\beta, \Sigma = \mathbb{E}[T(X)T(X)^\top] = \dots = \mathbb{E}[(AX)(AX)^\top] = A \mathbb{E}[X X^\top] A^\top = A \Sigma_\alpha A^\top$$

Lemma: $T(x) = A(x - m_\alpha) + m_\beta = \nabla \varphi(x)$, φ conv $\Rightarrow A^\top = A$ and $A \geq 0$

Cor: T is the OT $\Leftrightarrow \begin{cases} A \Sigma_\alpha A = \Sigma_\beta \\ A \geq 0 \end{cases}$ (Riccati equation)

RK: If $A > 0$, $\exists! \sqrt{A} \geq 0$ s.t. $\sqrt{A}^2 = A$

Proof: $A = U \text{diag}(\sigma^2) U^T$, $\sqrt{A} = U \text{diag}(\sigma) U^T$

Prop: $A \varepsilon_\alpha A^\top \varepsilon_\beta$ has a unique sol., $Q(\varepsilon_\alpha, \varepsilon_\beta)$.

Proof: $(\sum_{\alpha} A \varepsilon_\alpha)(\sum_{\alpha} A^\top \varepsilon_\alpha) = \sum_{\alpha} \varepsilon_\beta \sum_{\alpha} \varepsilon_\alpha \Leftrightarrow \sum_{\alpha} A \varepsilon_\alpha = \sqrt{\sum_{\alpha} \varepsilon_\beta \varepsilon_\alpha}$

(we are using ε has a density $\Rightarrow \det(\varepsilon_\alpha) \neq 0$).

$$\Leftrightarrow A = \sqrt{\sum_{\alpha} \varepsilon_\alpha}^{-1} \sqrt{\sum_{\alpha} \varepsilon_\beta \varepsilon_\alpha} \sqrt{\sum_{\alpha} \varepsilon_\alpha}^{-1} = Q(\varepsilon_\alpha, \varepsilon_\beta)$$

Def: If $AB = BA$, $G(A, B) := Q(A \cdot \varepsilon, B) = \sqrt{AB}$

Def: $W_2(\alpha, \beta) = \int \|x - T^*(x)\|^2 d\alpha(x)$, T^* Brenier map.

Prop: $W_2(\alpha, \beta) = \|m_\alpha - m_\beta\|_2^2 + \beta^2(\varepsilon_\alpha, \varepsilon_\beta)$ (see exercise sheet 1)

$$B(\varepsilon_\alpha, \varepsilon_\beta) := \sqrt{\varepsilon_\alpha + \varepsilon_\beta - \sqrt{\varepsilon_\alpha \varepsilon_\beta}}$$

Prop: if $\varepsilon_\alpha \varepsilon_\beta = \varepsilon_\beta \varepsilon_\alpha$, $B(\varepsilon_\alpha, \varepsilon_\beta) = \|\sqrt{\varepsilon_\alpha} - \sqrt{\varepsilon_\beta}\|_2 \geq \text{Hell}(G_\alpha, G_\beta)$,

$$\text{Hell}(G_\alpha, G_\beta) = \|G_\alpha - G_\beta\|_2$$

Kantorovich Formulation

cost of transport plan

Def: $P \in U(a, b)$ is a transport plan when $P \in \mathbb{R}_+^{n \times m}$ (positivity) and $\begin{cases} \forall i, \sum_j P_{ij} = a_i \\ \forall j, \sum_i P_{ij} = b_j \end{cases}$



$$U(a, b) := \{P \in \mathbb{R}_+^{n \times m} : P \mathbf{1}_m = a, P^T \mathbf{1}_n = b\}$$

c_{ij} is the cost for a unit amount of mass to be transported $x_i \rightarrow y_j$

$$(K) \min_P \left\{ \sum_{i,j} P_{ij} \cdot c_{ij} : \langle C, P \rangle = \langle a, P \mathbf{1}_m \rangle = a \right\}$$

Lemma: $U(a, b)$ is a non-empty compact convex set

Cor: (K) always admits a solution. (continuous function on a compact set)

This is the origin of linear programming \Rightarrow Kantorovich (theory, duality)

We study a particular case: $m=n$, $a=b=1$

$$(M) \min_{P \in \mathbb{R}_+^n} \sum_{i,j} c_{ij} P_{ij}, \quad P_n = \{P \in \mathbb{R}_+^{n \times n} : \forall i, j, P_{ij} \geq 0, \sum_j P_{ij} = 1\} \quad (\text{permutations})$$

$$(M) = \min_{P \in P_n} \langle P, C \rangle$$

$$B_n = \{P \in \mathbb{R}_+^{n \times n} : P \mathbf{1}_n = P^T \mathbf{1}_n = \mathbf{1}_n\} \Rightarrow (K) \min_{P \in B_n} \langle P, C \rangle$$

Theo.: Birkhoff - Von Neumann
 $(M) \subseteq (\mathcal{E})$, i.e. $\exists P$ solving (K) that is also a solution of (M) .

Def.: $P \in \text{Ext}^n(\mathcal{E}) \Leftrightarrow [P = \frac{Q+R}{2}, Q, R \in \mathcal{E} \Rightarrow Q=R=P]$ (vertices)

Theo.: If \mathcal{E} is compact, then $\text{Ext}^n(\mathcal{E}) \neq \emptyset$.

Theo.: (Krein - Milman)

\mathcal{E} compact and convex $\Rightarrow B = \text{ConvHull}(\text{Ext}^n(\mathcal{E}))$

Cor.: $\left\{ \min_{P \in \mathcal{E}} \langle c, P \rangle \right\} \cap \text{Ext}^n(\mathcal{E}) \neq \emptyset$

Proof: $S := \arg\min \{ \langle c, P \rangle \mid P \in \mathcal{E} \}$
 S compact and convex. $S \supset \text{Ext}^n(S) \in \text{Ext}^n(B)$
 $P \in \text{Ext}^n(S) \Rightarrow P \in \text{Ext}^n(\mathcal{E}) \quad P = \frac{Q+R}{2} \Rightarrow Q, R \in S$

Theo.: (B -VN)

$$\mathcal{J}_n^n = \text{Ext}^n(B_n)$$

Proof: i) $P \in \text{Ext}^n(B_n) \quad P_j = \frac{Q_j+R_j}{2}, Q_j, R_j \in [0, 1]^{n \times n} \rightarrow Q_{ij}, R_{ij} \in \{0, 1\}$, because $P_{ij} \in \{0, 1\}$
ii) $B_n \setminus P_n \supseteq B_n \setminus \text{Ext}^n(B_n)$

$P \in \text{Ext}^n(B_n), P \neq P_n$. Need to show $\exists Q \neq R \quad P = \frac{Q+R}{2}$

By graph theory arguments, there is at least one cycle in the transport graph.
 $i_1, j_1, \dots, i_k, j_k, \dots, i_1$

$$0 \leq \min_{ij} \{P_{ij}, 1 - P_{ij}\}$$

$$Q_{ij} = \begin{cases} P_{ij} + \epsilon & \text{if } (i, j) \in A \\ P_{ij} - \epsilon & \text{if } (i, j) \in B \\ P_{ij} & \text{otherwise} \end{cases}$$

$$A = \{(i_s, j_s)\}$$

$$B = \{(j_s, i_{s+1})\}$$

Lecture 3

Kantorovich: $\alpha = \sum_i \delta_{x_i}$, $\beta = \sum_j \delta_{y_j}$ $\rightarrow \min_{P \in \mathcal{P}^{\text{unif}}} \{ \langle \alpha, P \rangle : P \mathbb{I}_{m \times n}, P^\top \mathbb{1}_n = b \}$

We will now look at the generalised problem:

Def: $\alpha \in \mathcal{P}(X)$, $\beta \in \mathcal{P}(Y)$. A coupling is $\pi \in \mathcal{P}(X \times Y)$ having marginals α and β , i.e. $\forall f \in \mathcal{B}(X)$
 $\int f(x) d\pi(x, y) = \int f(x) d\alpha(x)$, $\forall g \in \mathcal{B}(Y)$, $\int g(y) d\pi(x, y) = \int g(y) d\beta(y)$

Def: $P_1 := (x, y) \in X \times Y \mapsto x \in X$, $P_2 := (x, y) \in X \times Y \mapsto y \in Y$
 $\pi_1 := (P_1) \# \pi$, $\pi_2 := (P_2) \# \pi$

Prop: Coupling $\Rightarrow \pi_1 = \alpha$, $\pi_2 = \beta$

Do it!

Prop: $\alpha \otimes \beta$ works

Prop: If marginals are discrete, $\begin{cases} \pi_1 = \beta \\ \pi_2 = \alpha \end{cases} \Leftrightarrow \pi = \sum_{ij} p_{ij} \delta_{(x_i, y_j)}$, $P \mathbb{I} = \alpha$, $P^\top \mathbb{1} = b$

$X = \bigcup_{j=1}^J J^2 j$ disjoint, $T: x \in J^2 j \rightarrow y_j$ piecewise constant.

Prop: iff $\alpha(J^2 j) = b_j$, $T \# \alpha = \beta$



Prop: Suppose $T \# \alpha = \beta$, then $\pi = \overbrace{(\text{Id}, T)}^Q \# \alpha$ is a valid coupling. $q: x \in X \mapsto (x, T(x)) \in X \times Y$

$$(K) \inf_{\pi \in \mathcal{P}(X \times Y)} \left\{ \int \int c(x, y) d\pi(x, y) : \pi_1 = \alpha, \pi_2 = \beta \right\}$$

RK: if (K) (discrete case) holds, $\pi^* := \sum_{ij} p_{ij}^* \delta_{(x_i, y_j)}$ and P^* is a solution to the discrete Kantorovich problem.

Prop: If α and β have compact support (i.e. X, Y compact), then (K) has a solution.

Proof: When we introduce topology:

For the weak* topology $\mathcal{P}(X \times Y)$ is compact, $\Pi \mapsto \int c d\Pi$ is continuous and $\alpha \otimes \beta$ is a valid coupling.

Prop: For a non compact space, $c(x, y) = d(x, y)^\rho$, $\int d(x, y)^\rho d\alpha(x) < +\infty$, $\int d\beta < +\infty$, then (K) has a solution.

Theo: (Brzennik)

$X = Y = \mathbb{R}^d$, $c(x, y) = \|x - y\|^2$, α has a density w.r.t. Lebesgue.

Ex! solution $T = \nabla q$ of Monge, $(\nabla q) \# \alpha = \beta$, $\pi = (\text{Id}, T) \# \alpha$ is the unique sol. of (K).

RK: $\pi = (\text{Id}, T) \# \alpha \Rightarrow \int \int c d\pi = \int c(x, T(x)) d\alpha(x)$.

RK: (K) $\hookrightarrow \inf_{(X, Y)} (E[c(X, Y)], \text{Law}(X) = \alpha, \text{Law}(Y) = \beta)$

Def: $X = Y$, $d(x, y) = d(x, y)^p$, $p \geq 1$.

Wasserstein distance: $W_p(\alpha, \beta) := \left(\inf_{\pi \in \mathcal{P}(X \times X)} \int \int d(x, y)^p d\pi(x, y) \right)^{1/p}$

Theo: $(\mathcal{P}(X), W_p)$ is a complete separable metric space.

Proof: $W_p(\alpha, \beta) = 0 \Leftrightarrow \alpha = \beta$. $(\alpha, \beta, \gamma) \in \mathcal{M}(X)^3$, $W_p(\alpha, \gamma) \leq W_p(\alpha, \beta) + W_p(\beta, \gamma)$

We look into the discrete case.

$$\alpha = \sum a_i \delta_{x_i}, \quad \beta = \sum b_j \delta_{y_j}, \quad \gamma = \sum c_k \delta_{z_k}.$$

$$\sum_{ij} P_{ij} \gamma_{ij}^p = 0 \Leftrightarrow P = \text{diag}(h). \quad P_1 = a \circ h, \quad P_2 = b \circ h \Rightarrow a = b$$

$$\underbrace{a \xrightarrow{P} b}_{=c} \xrightarrow{Q} c$$

$$\text{Lemma (Claim): } G_{ijk} = \frac{P_{ij} Q_{jk}}{b_j}, \quad \sum_i G_{ijk} = P_{ij}, \quad \sum_j G_{ijk} = Q_{jk} \quad \text{and} \quad \sum_i G_{ijk} a_i = c_j, \quad \sum_j G_{ijk} b_j = c_k, \quad \sum_k G_{ijk} = c_k$$

Def./Coro.: $P_{ik} = \sum_j G_{ijk}$ and we get $P_1 = a$, $P_2 = c$.

$$(\Leftrightarrow P = P \text{diag}\left(\frac{1}{b_j}\right) Q)$$

So: Take P, Q optimal $\rightarrow P$ coupling between a couple $(D_{\alpha\beta} = d(\alpha_x, \beta_y))$

$$W_p(\alpha, \beta)^p \leq \left(\sum_i D_{\alpha i}^p \frac{P_{ij} Q_{jk}}{b_j} \right)^{1/p} \leq \left(\sum_i (D_{\alpha i} \cdot D_{\beta i})^p \frac{P_{ij} Q_{jk}}{b_j} \right)^{1/p}.$$

$$\text{We use Markov's: } \left(\sum_i x_i (a_i - b_i) \right)^2 \leq \left(\sum_i x_i a_i^2 + \sum_i x_i b_i^2 \right).$$

$$\text{So } W_p(\alpha, \beta)^p \leq \left(\sum_i D_{\alpha i}^p G_{ijk} \right)^{1/p} + \left(\sum_i D_{\beta i}^p G_{ijk} \right)^{1/p} = W_p(\alpha, \beta) + W_p(\beta, \gamma).$$

$$\begin{aligned} & \sum_i D_{\alpha i}^p P_{ij} \\ & = W_p(\alpha, \beta) \end{aligned}$$

$$\text{Thm: } \Pi \in \mathcal{P}(K \times X), \quad \xi \in \mathcal{P}(K \times X) \Rightarrow \exists \varrho \in \mathcal{P}(X \times K \times X), \quad \varrho_{12} = (P_{xz})_* \varrho = \Pi, \quad \varrho_{23} = (P_{xz})_* \varrho = \xi$$

Def: (Total variation)

$$\|\alpha - \beta\|_{TV} = \|\alpha - \beta\|(X)$$

Introducing comparison:

$$(x_u) \rightarrow x. \quad \|\delta_{x_u} - \delta_x\|_{TV} \rightarrow 0, \quad W_p(\delta_{x_u}, \delta_x) = d(x_u, x) \rightarrow 0$$

Def: Convergence in TV is strong convergence.

Def: $\alpha_k \xrightarrow{*} \alpha \Leftrightarrow \forall f \in C_b(X), \int f(x) d\alpha_k(x) \rightarrow \int f(x) d\alpha(x)$

Prop: $\|\alpha\|_{TV} = \sup_{\|f\|_{\infty} \leq 1} \int f(x) d\alpha(x)$

Theo: Wasserstein is cv in law.

If X compact, $\alpha_k \xrightarrow{*} \alpha \Leftrightarrow W_p(\alpha_k, \alpha) \xrightarrow{k \rightarrow \infty} 0$

In general: $W_p(\alpha_n, \alpha) \rightarrow 0 \Leftrightarrow \int \alpha_n \xrightarrow{*} \alpha$

$$\int d(x_n, x)^p d\alpha_n \rightarrow \int d(x_n, x)^p d\alpha$$

Theo: If $\mathbb{E}[|X|^3] < \infty$, then $W_p(\alpha_n, \mu) \leq \frac{C}{\sqrt{n}}$ (α_n the law of $\frac{X_n - \mu}{\sigma_n}$)

Lecture 4

We regularise the problem:

$$\min_{P \in \mathbb{R}^{n \times m}} \{ \langle c, P \rangle + \epsilon H(P) : P \mathbf{1} = a, P^T \mathbf{1} = b \} \quad (\text{Schrödinger problem})$$

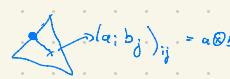
Def: Shannon entropy: $H(P) := \sum_{ij} P_{ij} \log(P_{ij})$

RK: As ϵ increases our graph becomes more and more connected.

Prop: P is the solution of (Schw) iff. $\begin{cases} P \mathbf{1} = a \\ P^T \mathbf{1} = b \\ \exists u \in \mathbb{R}_+^n, v \in \mathbb{R}_+^m, P = \text{diag}(u) K \text{diag}(v) \end{cases}$

Def: $K = \exp\left(-\frac{c}{\epsilon}\right) = \left(\exp\left(-\frac{c_{ij}}{\epsilon}\right)\right)_{ij}$

RK: For W_2 , K is the Gaussian kernel.

Proof: Step 1: $P_{ij} > 0 \quad \forall i, j \in \text{supp } P_{ij} = 0$ 

$$P_{ij}^+ = (I - A) P_{ij}^- + (a_i b_j)^T, \quad f(+): \langle c, P^+ \rangle + \epsilon H(P^+)$$

Prop: $f(0) = -\infty$  This implies that there is a better solution than the one with $P_{ij}^+ = 0$.

Now: $\min_{P \geq 0} \{ \langle c, P^+ \rangle \in H(P) : P \mathbf{1} = a, P^T \mathbf{1} = b \}$ 

$$L(P, f, g) = \langle c, P \rangle + \epsilon H(P) + \underbrace{\langle a - P \mathbf{1}, f \rangle}_{\langle a, f \rangle - \langle P, f \mathbf{1}^T \rangle} + \underbrace{\langle b - P^T \mathbf{1}, g \rangle}_{\langle b, g \rangle - \langle P^T \mathbf{1}, g^T \rangle}$$

$$\frac{\partial L}{\partial P} = 0 = c + \epsilon(\log(P) + I) - f \mathbf{1}^T - g \mathbf{1}^T = 0 \quad \text{and get the formula:} \\ P_{ij} = \exp\left(-\frac{c_{ij}}{\epsilon}\right) \cdot \exp(f_i/\epsilon) \cdot \exp(g_j/\epsilon) \cdot e$$

\odot is element-wise multiplication

$$\begin{cases} (P) \quad u \odot (K \alpha) = a \rightsquigarrow u = \frac{a}{K \alpha} & (\text{element-wise division}) \\ (C) \quad \alpha \odot (K^T u) = b \rightsquigarrow \alpha = \frac{b}{K^T u} \end{cases}$$

Many names for the algorithm: Sinkhorn, Iterative scaling, TPFP

$$u_{++} = 1 \rightarrow u_{++} = \frac{a}{K \alpha_{++}} \rightarrow \alpha_{++} = \frac{b}{K^T u_{++}} \boxed{\text{O}(T \cdot n^2)} \quad \# \text{ of iterations}$$

Convergence? (spoiler: yes)

• Prop.: to compute the Wasserstein distance to a distribution S , run S (up to log (d)) and run $\frac{1}{\delta^2} \approx \frac{1}{S}$ iterations of Sinkhorn. ($\| \langle Q, P_S \rangle - \langle \cdot, P^* \rangle \|_2^2 \leq \delta$).

• Def.: $KL(P||Q) = \sum_i p_i \log \left(\frac{p_i}{q_i} \right)$, $H(P) = KL(P||P)$

• Prop.: $KL(P||Q) \geq 0$, $KL(P||Q) = 0 \Rightarrow P = Q$

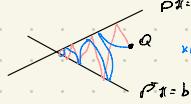
• Prop.: $P \# a = a$, $P^\top \# b = b$, $KL(P \# a \otimes b) = KL(P \# a \otimes b^\top) = KL(a \# a^\top) + KL(b \# b^\top)$
Therefore, taking $a = z$, $b = z$, $KL(P \# a \otimes b) \geq H(P) + \text{cst}$

(Sch) $\in \min_{P \geq 0} (\langle \cdot, P \rangle + \epsilon KL(P \# a \otimes b)) : P \# a = a, P^\top \# b = b$

• Prop.: $P \xrightarrow{\epsilon \rightarrow 0} a \otimes b$

• Prop.: $\langle \cdot, P \rangle + \epsilon KL(P \# a \otimes b) = \epsilon KL(P) e^{-\epsilon H(P \# a \otimes b)} + \text{cst}$

This is all a story of projection on convex sets.



$\text{Proj}_{E_\alpha}^{KL}(P) = \inf_{P \in E_\alpha} KL(P||P')$, $P_0 = Q$, $P_{t+\frac{1}{2}} = \text{Proj}_{E_\alpha}^{KL}(P_t)$, $P_{t+1} = \text{Proj}_{E_\alpha}^{KL}(P_{t+\frac{1}{2}})$

• Theo.: (Bregman)

$P_t \rightarrow E_\alpha \cap E_\beta$ (the constraints). If E_α and E_β are affine, $P_t \rightarrow \text{Proj}_{E_\alpha \cap E_\beta}(Q)$

• Prop.: This is Sinkhorn.

$\begin{cases} E \text{ strictly convex}, \\ B_\varphi(\tau(y)) = \varphi(\tau) - \varphi(y) - \nabla \varphi(\tau)^\top y, \nabla \varphi(\tau) > 0 \\ \text{we can build many things from that: } q(s) = \frac{\|s\|^2}{2} \Rightarrow B_\varphi(\tau(y)) = \frac{1}{2} \|s-y\|^2, \varphi(s) = f(s) \Rightarrow B_\varphi = KL \end{cases}$

• Def.: $E := \mathbb{R}_+^d / \sim$, $u \sim v \Leftrightarrow \exists \lambda \geq 0: u = \lambda v$
 $d_H(u, v) = \|\log(\frac{u}{v})\|_1$, $\|z\|_1 = \max(z) - \min(z)$

• Theo.: (E, d_H) is a complete metric space.

• Theo.: (Birkhoff contraction)

$K: E \rightarrow E$ linear. $\exists m \in \mathbb{R}$, $d_H(Ku, Ku') \leq m d_H(u, u')$

• Prop.: $d_H(\frac{u}{m}, \frac{u'}{m}) = d_H(u, u')$

• Cor.: $d_H(s(u), s(u')) \leq m^2 d_H(u, u')$ ($s(u)$ means one iteration of Sinkhorn)

• Prop.: $m = \exp(-\frac{\gamma \|C\|_1}{\epsilon}) \xrightarrow{\epsilon \rightarrow 0} 1$

$$\alpha, \beta \in M^+(X), M^+(Y), \inf_{\pi} \left\{ \iint_{X \times Y} c d\pi + \epsilon KL(\pi \| \alpha \otimes \beta) \mid \pi_x = \alpha, \pi_y = \beta \right\}$$

Theo: If c continuous and X, Y compact, then (Schmid) has a unique solution π_* .

$$\left[\frac{\partial \pi}{\partial \alpha \partial \beta}(x, y) = u(x) \exp\left(-\frac{c(x, y)}{\epsilon}\right) v(y) \right]$$

Def: $X, Y \sim \pi$, $I(X, Y) := K\{H | \pi_* \otimes \pi_*\}$

Prop: $I(X, Y) = 0$ iff. $X \perp\!\!\!\perp Y$

$$(sch) \quad \inf_{(X, Y)} \left\{ E[c(X, Y)] + \epsilon I(X, Y) : X \sim \alpha, Y \sim \beta \right\}$$

Lecture 5

Kernels: $X = Y$, $K(x, y) \in \mathbb{R}$ that is positive definite.

Def.: $K(x, y)$ is universal iff: $\text{span}\{x \mapsto K(x, y) : y \in X\}$ is dense in $C(X, \mathbb{R})$
 $\forall X \subseteq \mathbb{X}$ compact, for uniform convergence on B .

Ex.: $K(x, y) = \exp(-\frac{\|x-y\|^2}{\sigma^2})$

RK: PD Kernel $\Leftrightarrow K(x, y) = \langle \varphi(x), \varphi(y) \rangle_{\mathcal{H}}$

Def: Maximum mean discrepancy (MMD)

$$\|\alpha - \beta\|_K^2 = \iint K(x, y) d(\alpha - \beta)(x) d(\alpha - \beta)(y)$$

RK: $K(x, y) = h(x-y) \rightarrow \int h(x-y) d\gamma(y) = (h * \gamma)(x)$

$$\iint K d\gamma \otimes \gamma = \langle h * \gamma, \gamma \rangle$$

Theo: Def. 4

If K is conditionally pd. and universal, then $\|\cdot\|_K$ is a norm on $P(X)$.

$$\alpha_s \rightarrow \alpha \Leftrightarrow \|\alpha - \alpha_s\|_K \rightarrow 0$$

Discrete case:

$$\alpha = \frac{1}{n} \sum_i \delta_{x_i}, \quad \beta = \frac{1}{m} \sum_j \delta_{y_j}, \quad \|\alpha - \beta\|_K^2 = \|\alpha\|_K^2 + \|\beta\|_K^2 - 2 \langle \alpha, \beta \rangle_K = \frac{1}{N^2} \sum_{i,j} K(x_i, y_j)$$

$$+ \frac{1}{m^2} \sum_{j,j'} K(y_j, y_{j'}) - \frac{2}{nm} \sum_{i,j} K(x_i, y_j)$$

$$W_c^\epsilon(\alpha, \beta) = \min_{\substack{\pi_{\epsilon} \text{ s.t.} \\ \pi_{\epsilon} = \alpha}} \left(\iint c d\pi_{\epsilon} + \epsilon K_L(\pi_{\epsilon} \| \alpha \otimes \beta) \right) \xrightarrow{\epsilon \rightarrow 0} W_c(\alpha, \beta)$$

$$\Delta W_c^\epsilon(\alpha, \alpha) > 0.$$

$$\pi_\epsilon = \text{argmin } W_c^\epsilon(\alpha, \beta)$$

$$\text{Prop: } \pi_\epsilon \xrightarrow{\epsilon \rightarrow 0} \alpha \otimes \beta, \quad W_c^\epsilon(\alpha, \beta) \xrightarrow{\epsilon \rightarrow 0} \iint c(x, y) d\alpha \otimes \beta(x, y).$$

$$\text{Def: } K_c(x, y) = -c(x, y)$$

Def: Sinkhorn divergence (Thibault Sejourne)

$$\bar{W}_c^\epsilon(\alpha, \beta) = W_c^\epsilon(\alpha, \beta) - \frac{1}{\epsilon} W_c^\epsilon(\alpha, \alpha) - \frac{1}{\epsilon} W_c^\epsilon(\beta, \beta)$$

$$\text{Prop: } \bar{W}_c^\epsilon(\alpha, \beta) \xrightarrow{\epsilon \rightarrow 0} W_c(\alpha, \beta) \geq 0$$

$$\bar{W}_c^\epsilon(\alpha, \beta) \xrightarrow{\epsilon \rightarrow 0} -\langle \alpha, \beta \rangle_K + \frac{1}{2} \|\alpha\|_K^2 + \frac{1}{2} \|\beta\|_K^2 = \frac{1}{2} \|\alpha - \beta\|_K^2$$

Theo: If $\exp(-\frac{c(x,y)}{\epsilon})$ is a valid kernel, then $\bar{w}_c \geq 0$, $\bar{w}_c^T(\alpha, \beta) = 0 \Rightarrow \alpha = \beta$,
 $\alpha_s \in \alpha \Leftrightarrow \bar{w}_c^T(\alpha_s, \alpha_s) \rightarrow 0$.

Theo: $d = w_p$, $A_n \sim \frac{1}{n} \delta_{xy}$, $\Delta_n = E[d(\alpha, \beta)] - d(\hat{\alpha}, \hat{\beta})$

where $\hat{\alpha}, \hat{\beta}$ are random measures sampling from α, β .

If $d = \text{MMD}$, $\Delta_n \leq \frac{C_d}{\sqrt{n}}$

If $d = w_c^T$, $\Delta_n \leq C_d(1 + \epsilon^{-\frac{d}{2}}) \cdot \frac{1}{\sqrt{n}}$

Theo: if $c(x, y)$ is Lips., X, Y compact, then $\min_{x \in X} \int c(x, y) d\pi(x) = \sup_{\substack{f \in \mathcal{C}(X) \\ g \in \mathcal{C}(Y)}} \left\{ \int f(x) d\pi(x) + \int g(y) dB(y) \right\}$: $\forall x, y, f(x) + g(y) \leq c(x, y)$

Def: $g \in \mathcal{C}(Y) \rightarrow g^c(x) = \min_{y \in Y} (c(x, y) - g(y)) \in \mathcal{C}(X)$ from \mathcal{C}(X) to \mathcal{C}(Y)

$\therefore \sup_{g \in \mathcal{C}(Y)} (\int d\pi : f(x), g(y) \leq c(x, y)) = \sup_g (\int g(y) dB(y)) : g(y) \in \bar{f}^c(y)$

Prop: $g = \bar{f}^c$ β -a.e.

Summary: f fixed $\rightarrow f^c$ soln. g fixed $\rightarrow g^c$ solution.

Def: Legendre transform

$$f^*(y) = \sup_x (x \cdot y - f(x))$$

Def: c -concave functions := $\{f : \exists g, f = g^c\}$ ($c(x, y) = L(x, y) \rightarrow$ concave = c -concave)

Prop: f c -concave $\Rightarrow f^c = f$

Lecture 6

1 mysterious have because of an accident.

Def.: $\delta f(\alpha) : \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(\alpha + \varepsilon e_i) = f(\alpha) + \varepsilon \underbrace{\langle e_i, \delta f(\alpha) \rangle}_{\int \delta f(\alpha)(x) d\mu(x)} + o(\varepsilon)$$

$$\frac{\partial \alpha}{\partial t} = -\delta f(\alpha)$$

$$\text{Ex.: } f(\alpha) = \int h(x) d\mu(x), \quad \delta f(\alpha) = h \Rightarrow \frac{\partial \alpha}{\partial t} = -h$$

$$f(\alpha) = \int |\nabla \alpha(x)|^2 d\mu = -\Delta \alpha, \quad \Rightarrow \delta f(\alpha) = -\Delta \alpha$$

$\alpha_{t+\tau}$: classical flow of Dirichlet is the heat equation $\frac{\partial \alpha}{\partial t} = \Delta \alpha$.

$$\alpha_{t+\tau} = \min_{\alpha} \frac{1}{2\tau} W_2^2(\alpha_t, \alpha) + f(\alpha)$$

$\tau \rightarrow 0$ gives the Wasserstein gradient flow of f .

$$\alpha = \delta_x, \quad F(\tau) = f(\delta_x)$$

$$\tau_{t+\tau} = \min_{\alpha} \frac{1}{2\tau} \|x - x_{t+\tau}\|^2 \in F(\tau) \quad (\text{Implicit Euler})$$

The opt. cond. is $(\tau_{t+\tau} - x_{t+}) + \tau \nabla F(x_{t+}) = 0$

$$(IE) \quad x_{t+\tau} = (\text{Id} + \tau \nabla F)^{-1}(x_t) \quad || \quad (EE) \quad \tau_{t+\tau} = (\text{Id} - \tau \nabla F)(x_t)$$

Theorem:

Under appropriate hypothesis, the W_2 -GF (Gradient Flow) is the solution of:

$$\frac{\partial \alpha_t}{\partial t} + \text{div}(\alpha_t (-\nabla_W f(e_t))) = 0$$

$$\text{where } \nabla_W f(\alpha) = \nabla_{\mathbb{R}^n} [\delta f(\alpha)]$$

$$\text{Linear: } f: f(\alpha) = \int h(x) d\mu(x), \quad \delta f(\alpha) = h, \quad \nabla_W f(\alpha)(x) = \nabla h(x)$$

$$\frac{\partial \alpha}{\partial t} - \text{div}(\nabla h \alpha) = 0 \Rightarrow \dot{\alpha} = -\nabla h \alpha$$

$$\text{Shannon: } f(\alpha) = \log \left(\frac{d\alpha(x)}{dx} \right) d\mu(x), \quad \delta f(\alpha) = \log(\alpha(x)) + 1, \quad \nabla_{\mathbb{R}^n} \delta f(\alpha) = \frac{\partial \alpha}{\partial x}(x)$$

$$\boxed{\frac{\partial \alpha}{\partial t} = \Delta \alpha} \quad (\text{heat equation})$$

$$\dots \approx \inf_{\alpha} \left\{ f(\alpha_*) + \epsilon \right\} \left[\frac{1}{2} \alpha^T + \epsilon \alpha, \nabla f(\alpha_*) \right] d\alpha_* + o(\epsilon) \}$$

More: $f(\alpha) = \int g\left(\frac{d\alpha}{dx}\right) dx$, $\nabla f(\alpha) = g^T(\alpha(x))$, $\nabla_\alpha f(\alpha) = g^T(\alpha(x)) - \nabla \alpha(x)$

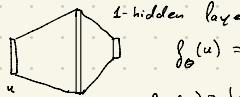
$$\frac{\partial \alpha}{\partial t} + \operatorname{div}(\nabla \alpha(t)) \cdot h(\alpha) = 0, \quad \text{where } h(s) = g^T(s) \cdot s$$

Def: f is geodesically convex when $\forall \alpha, \beta, \gamma \in \mathbb{R}^d$, $\Rightarrow f(\alpha + \beta) \leq f(\alpha) + f(\beta)$ is convex

Theo: $\int h d\alpha$ good convex as h convex
 $\rightarrow \int \alpha^T \log \alpha$ is good convex

Theo: (McLennan)

$f(\alpha) = \int g\left(\frac{d\alpha}{dx}\right)$ is good convex if $g(0)=0$, g convex \nearrow , $g(s^*) s^*$ is ∇ in \mathbb{R}^d .



$$f_\theta(u) = \frac{1}{n} \sum_{i=1}^n q_i G(u; \omega_i, \theta) = A\phi(\theta), \quad \theta_i = (\omega_i, b_i), \quad \alpha = \frac{1}{n} \sum_{i=1}^n \delta_{\theta_i},$$

$$(*) f_\theta(\alpha) = \int \phi(\theta, \alpha) d\alpha(\theta), \quad \phi(\theta, \alpha) = q(G(\alpha; \omega, \theta))$$

Prop: $E(\alpha) = \text{cost.} + \int b(\theta) d\alpha(\theta) + \int k(\theta, \alpha) d\alpha(\theta) + \int \ell(\theta, \alpha) d\alpha(\theta)$, $b(\theta) = \frac{1}{N} \sum_k y_k - \phi(\theta; u_k)$
 $k(\theta, \alpha) = \frac{1}{N} \sum_k q(\theta; u_k) \phi(\theta; u_k)$ (we look for $\min_{\alpha} E(\alpha)$)

Theo:

If α_{+0} has enough particles (in range) and $\alpha_+ \rightarrow \alpha_\infty$, then α_∞ is a global min.